

Changes in Neural Processing After Multimodal Speech-Gesture Training in Patients With Schizophrenia

Lydia Riedl^{1,2}, Arne Nagels³, Gebhard Sammer⁴, Momoko Choudhury¹, Annika Nonnenmann¹, Anne Sütterlin¹, Chiara Feise¹, Maxi Haslach¹, Florian Bitsch¹, Hilde Kuehne⁵, Thomas S. Hartmann^{1,2}, Nina Shvetsova⁵, Benjamin Straube¹

¹Translational Neuroimaging Lab, Department of Psychiatry and Psychotherapy, Philipps-University Marburg, Germany; ²Center for Mind, Brain and Behavior (CMBB), Philipps-University Marburg and Justus Liebig University Giessen, Germany; ³Department of English and Linguistics, Johannes-Gutenberg-University Mainz, Germany; ⁴Department of Psychiatry and Psychotherapy, Justus-Liebig-University Giessen, Germany; ⁵CVAI Group, Goethe-University Frankfurt, Germany

Background

Dysfunctional social communication in patients with schizophrenia spectrum disorder (SSD):

- Interpreting abstract (abs) vs. concrete (con) speech
- Integrating speech and gesture (SG)
- aberrant activation in left inferior frontal gyrus (IFG) in superior/middle temporal regions (STS, MTG)¹

Objectives:

Improving integration of gesture in abstract speech context through a novel multimodal speech-gesture training (MSG):

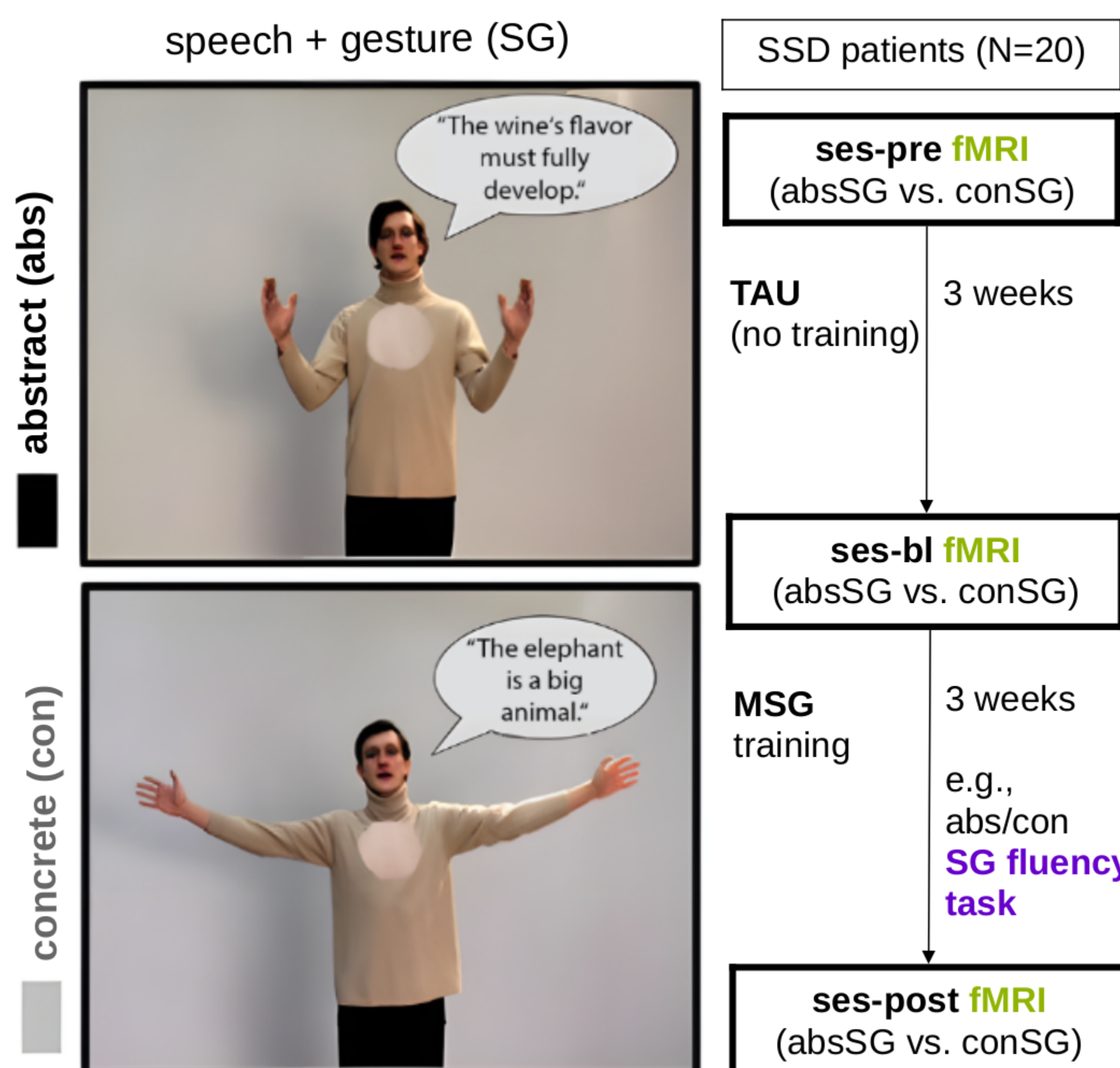
- 1) In behavioral performance
- 2) In neural processing
- 3) Explore associations of behavioral + neural changes

Methods

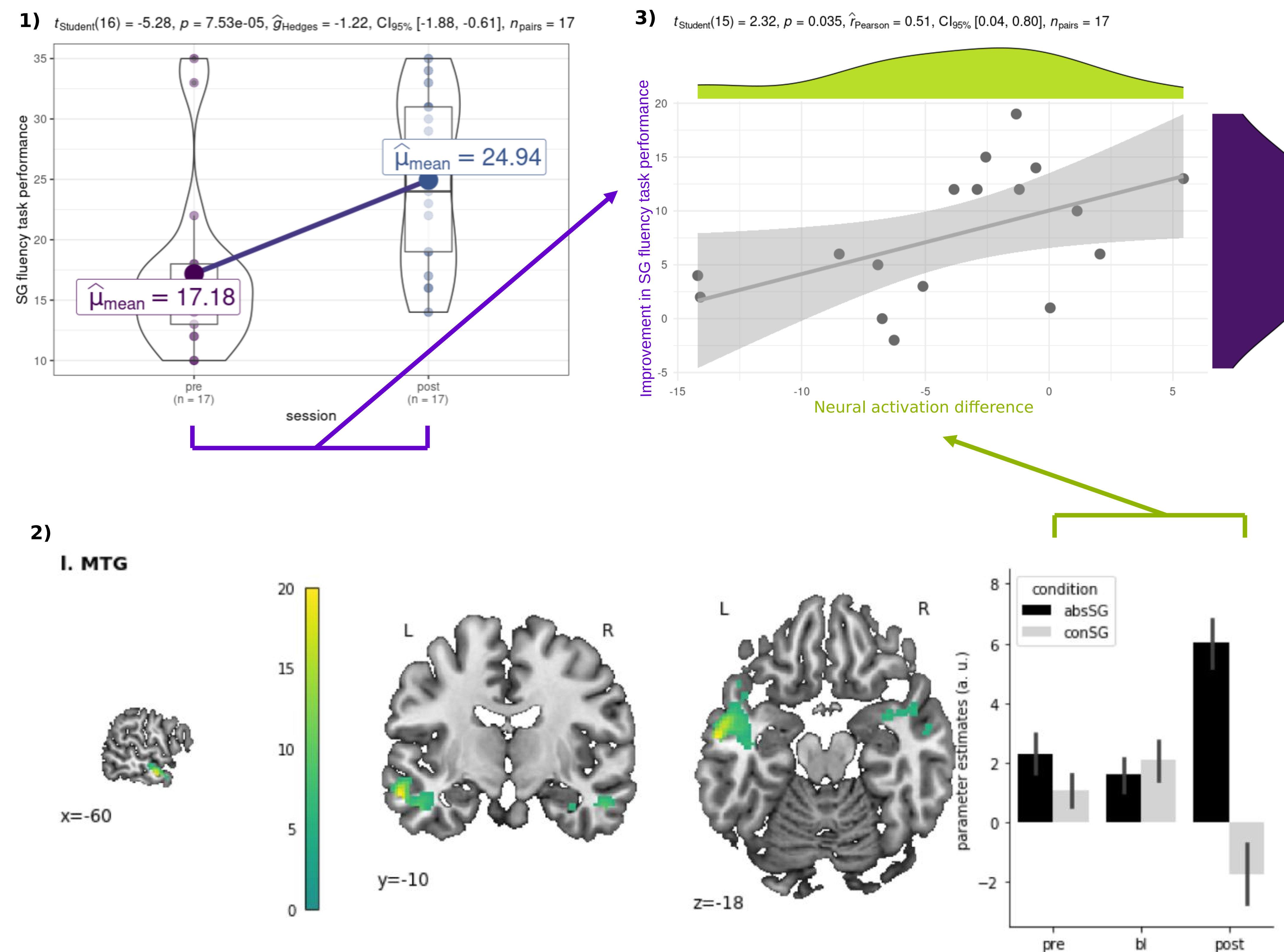
fMRI data acquisition: 3 Tesla scanner

Analysis: fMRIPrep² + SPM12; ggstatsplot³ (RStudio)

Contrast of interest: sessions X abstractness in SG



Results



1. MSG training session (first versus last training session) comparison of the patients' performance in speech-gesture (SG) fluency training task. The values represent the sum of correctly produced speech-gesture pairs for each of the three semantic categories per training session. Patients were asked to produce as many speech-gesture pairs as possible in one minute per semantic field.

2. Neural activation sessions X abstractness contrast. Activation clusters were thresholded at $p < 0.01$, with a minimum cluster size of 221 voxels, cluster level corrected at $p > 0.05$. Error bars indicate the standard error of the mean (s.e.m.). Crosshair point at the left MTG [$x = -60, y = -10, z = -18$]. MTG = middle temporal gyrus; L = left; R = right.

3. Correlation of the patients' improvement during training (last - first training session) in the SG fluency task and the difference (ses-post - ses-pre) in neural activation (sessions X abstractness contrast) in the conSG condition.

Discussion

- 1) **improvement in behavioral performance** (e.g., SG fluency training task): possible modification of dysfunctional social-communicative skills in SSD
 - 2) **sessions X abstractness interaction:** → possible training-specific effects
 - 3) **Correlation of 1) and 2):** → Possible association of MSG-related modification in neural activation
- Self-report measures and ratings from relatives confirm MSG-related changes.

Outlook

Investigation of the underlying mechanisms of changes in neural processing after MSG training:

- analyses of audio-visual concept learning in the Audio-Video Language Network (AVLnet)⁴, a self-supervised network that learns a shared audio-visual embedding space directly from raw video inputs.
- comparing AVLnet's learned representations with the MSG training induced changes
- implications for future therapy approaches?

References:

1. Straube B, Green A, Sass K, Kirner-Veselinovic A, Kircher T (2013): Neural integration of speech and gesture in schizophrenia: Evidence for differential processing of metaphoric gestures. Hum Brain Mapp 34:1696–1712.
2. Esteban O, Markiewicz CJ, Blair RW, Moodie CA, Isik AI, Erramuzpe A, Kent JD, Goncalves M, DuPre E, Snyder M, Oya H, Ghosh SS, Wright J, Durnez J, Poldrack RA, Gorgolewski KJ (2018): fMRIPrep: a robust preprocessing pipeline for functional MRI. Nature Methods 16:111–116.
3. Patil I (2021): Visualizations with statistical details: The "ggstatsplot" approach. Journal of Open Source Software 6:3167.
4. Rouditchenko A, Boggust A, Harwath D, Chen B, Joshi D, Thomas S, Audhkhasi K, Kuehne H, Panda R, Feris R, Kingsbury B, Picheny M, Torralba A, Glass J (2021): AVLnet: Learning Audio-Visual Language Representations from Instructional Videos. arXiv:200609199 [cs, eess]. <http://arxiv.org/abs/2006.09199>.

Acknowledgements:

The project is funded by the von behring|röntgen|foundation and supported by "The Adaptive Mind", funded by the Excellence Program of the Hessian Ministry of Higher Education, Science, Research and Art.

This study is preregistered (DRKS00015118).

Design and Methods are in detail described [here](#). Please find our public project incl. material and first data also on [OSF](#).

Contact: Riedl@staff.uni-marburg.de

 @RiedlLydia